# `EE-Tuning`: An Economical yet Scalable Solution for Tuning Early-Exit Large Language Models

Xuchen Pan*, Yanxi Chen*, Yaliang Li, Bolin Ding, Jingren Zhou

{panxuchen.pxc, chenyanxi.cyx, yaliang.li, bolin.ding, jingren.zhou}@alibaba-inc.com

Alibaba Group

**Abstract**

This work introduces `EE-Tuning`, a lightweight and economical solution to training/tuning *early-exit* large language models (LLMs). In contrast to the common approach of full-parameter pre-training, `EE-Tuning` augments any pre-trained (and possibly fine-tuned) standard LLM with additional early-exit layers that are tuned in a parameter-efficient manner, which requires significantly less computational resources and training data. Our implementation of `EE-Tuning` achieves outstanding training efficiency via extensive performance optimizations, as well as scalability due to its full compatibility with 3D parallelism. Results of systematic experiments validate the efficacy of `EE-Tuning`, confirming that effective early-exit LLM inference can be achieved with a limited training budget. In hope of making early-exit LLMs accessible to the community, we release the source code of our implementation of `EE-Tuning` at https://github.com/pan-x-c/EE-LLM.

## Contents

---

*Co-first authors.

# 1   Introduction

Transformer-based large language models (LLMs) have achieved extraordinary performance on various language tasks [51, 4, 32, 48, 49, 7]. Meanwhile, these models incur high costs and latency during the inference phase, due to their increasingly large sizes. *Early exiting* has proven to be a simple yet effective technique for accelerating inference of LLMs and other deep neural networks. In this approach, early-exit layers are attached to the original deep neural network, which can convert intermediate hidden states into early-exit output. During inference, the model can adaptively select one early exit to generate the output for each input sample, skipping the forward computation of the remaining layers of the network. Early exiting has found success in natural language processing [13, 18, 57, 41, 29, 11, 52, 27, 40, 53, 54, 19], computer vision [33, 47, 22, 21], and many other areas [38, 26, 14, 9].

This work considers token-wise early exiting for generative LLMs and autoregressive natural language generation [40, 8, 2, 50, 12, 6]. While the majority of prior works in this area have focused on designing early-exit *inference* mechanisms, we instead focus on how to *train* an early-exit LLM in the first place. The standard and straightforward method, adopted in most prior works on early exiting, is to jointly train all model parameters (including the network backbone and early-exit layers) from scratch, by minimizing a weighted sum of training losses from early and final exits. Recent work has made this approach compatible with massive 3D parallelism, thereby scaling up early-exit LLMs to sizes as large as any standard LLM that can possibly be trained with state-of-the-art LLM frameworks [43, 31, 6]. The obvious issue with this approach is its excessively high costs and complexity. Indeed, the massive amount of computational resources required to train an LLM with billions of parameters is simply inaccessible to most members of the community. Oftentimes in practice, one has access to the weights of an existing pre-trained (and possibly fine-tuned) standard LLM that might be private or open-source, and wonder if it is possible to train an early-exit LLM by leveraging such information rather than from scratch.

All these motivate us to *convert an existing generative LLM to an early-exit one*, in a way that

- requires minimum computational resources;

- leads to satisfactory inference acceleration; and

- preserves the full capability of the original LLM.

**Main contributions.**   This work introduces `EE-Tuning`, a principled and lightweight approach of transforming a pre-trained (and possibly fine-tuned) LLM into an early-exit one, which satisfies all the above requirements. At the core of `EE-Tuning` is an intuitive and practical two-stage procedure:

1. Take a pre-trained standard LLM as input, and augment its architecture with early-exit layers, whose parameters are initialized properly;

2. Tune the early-exit layers via backpropagation of certain training losses in a parameter-efficient manner, with modules of the original standard LLM frozen.

See Figure 1 for a visualization of our method, whose details will be elucidated in Section 2. Our implementation is based on the recently proposed EE-LLM framework [6], complementing the latter with an alternative solution to training early-exit LLMs that is both *accessible* and *scalable*, thanks to its low computational complexity and full compatibility with 3D parallelism. In other words, any LLM developer, with access to either one GPU or a cluster with thousands of GPUs, will find `EE-Tuning` a useful and practical tool for studying and applying early exiting. Our implementation also includes support for various configurations and other favorable features, which further makes it more convenient to use.[1]

The efficacy of `EE-Tuning` is validated via extensive and systematic experiments for models with up to 70 billion (70B) parameters, an unprecedented scale for early-exit LLMs. More specifically, a pre-trained LLM can quickly acquire the ability of early exiting via the tuning process with fast and stable convergence, which takes less than 1/1000 of the GPU hours and training data used in its pre-training stage, and requires only one or a few GPUs. Meanwhile, the converted model can achieve $1.2\times$ to $1.6\times$ speedup on various

---

[1]Henceforth, we let `EE-Tuning` refer to the proposed two-stage method or our implementation of it, depending on the context.
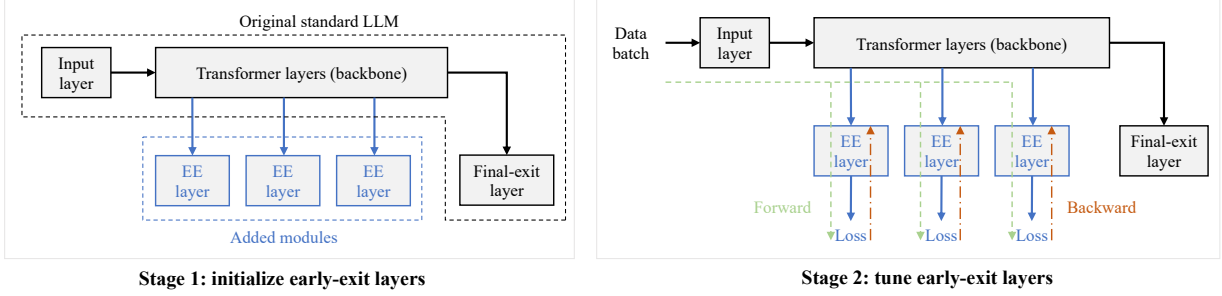
Figure 1: An outline of `EE-Tuning`, the proposed two-stage procedure that converts a pre-trained standard LLM into a well-trained early-exit LLM.

downstream tasks through early exiting while maintaining comparable or even better benchmark scores, or higher speedup if slight degeneration of output quality is acceptable. We thoroughly investigate the effects of various design choices and provide practical guidelines for maximizing the performance of `EE-Tuning`. The source code of `EE-Tuning` is available at https://github.com/pan-x-c/EE-LLM.

**Related works.** The idea of `EE-Tuning`, i.e. augmenting a pre-trained neural network with early-exit layers that are tuned in a parameter-efficient manner, is not completely new. This strategy has been adopted in some prior works for model architectures tailored to classification tasks, e.g. the encoder-only BERT model [52, 29, 19] or others [33, 22, 3, 9]. However, there is no guarantee that results and conclusions from these works can safely transfer to the case of *decoder-only* Transformers tailored to *autoregressive sequence generation*, which is the focus of our work. Another recent work [50] proposed to initialize the model parameters of early-exit LLMs with pre-trained standard LLMs, but followed by full-parameter training. Closest to our setting and training methodology is the recent work [12], which investigated generative LLMs of sizes up to 355M, and only considered linear exit heads that are randomly initialized. Moreover, that work proposed to use multiple exits for improving the final output of full-model inference, rather than accelerating inference via early exiting. In contrast, our implementation of the proposed `EE-Tuning` method (1) is unified and systematic, with support for a wide range of configurations; (2) is scalable, thanks to its full compatibility with 3D parallelism; and (3) has proven via extensive experiments to return early-exit LLMs that achieve outstanding acceleration during autoregressive inference.

## 2  Methodology

This section elaborates on our methodology and implementation of obtaining a well-trained early-exit LLM via `EE-Tuning`, the two-stage procedure visualized in Figure 1.

**Preliminaries.** Modern LLMs are mostly based on the Transformer architecture [51]. We focus on the decoder-only generative pre-training (GPT) Transformer architecture [35, 36], although many techniques presented in this work can be generalized to broader settings. As visualized in Figure 1, a GPT Transformer is composed of an initial layer for input processing, a stack of Transformer layers as the backbone, and an output layer that converts the final hidden states to logits on the vocabulary, which can be used for generating new tokens. Each Transformer layer consists of an attention module and a multi-layer perceptron (MLP), with layer normalization [1] and residual connections [15] applied in multiple places. The final output layer includes an optional layer normalization module, followed by a large output embedding matrix.

A GPT Transformer can be trained in an unsupervised manner, by optimizing the language modeling loss on unlabeled corpus. More specifically, given model parameters $\boldsymbol{\theta}$ and a sequence of tokens $\boldsymbol{x} = (x_1, x_2, \ldots, x_T)$, the autoregressive language modeling loss, i.e. negative log-likelihood of next-token prediction, is defined as $\mathcal{L}(\boldsymbol{x}; \boldsymbol{\theta}) \coloneqq -\log \mathbb{P}(\boldsymbol{x}; \boldsymbol{\theta}) = -\sum_{t \in [T]} \log \mathbb{P}(x_t | x_1, \ldots, x_{t-1}; \boldsymbol{\theta})$.
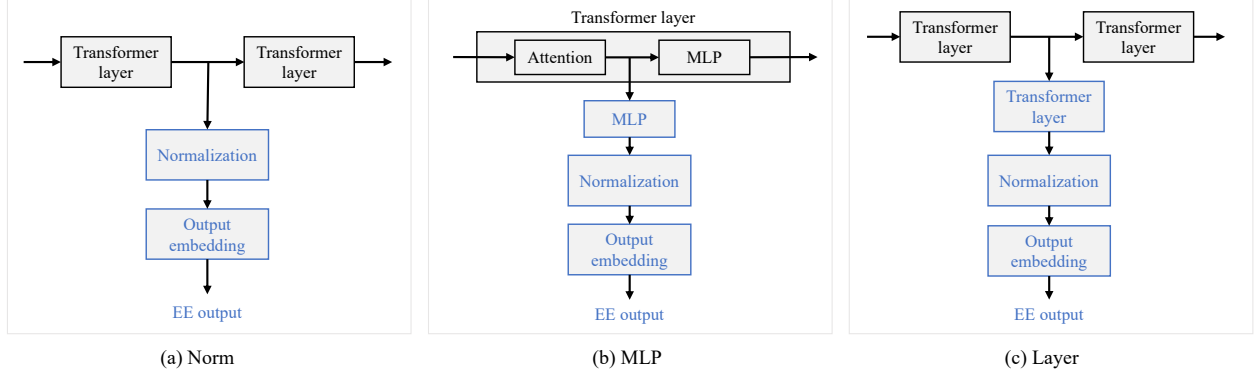
Figure 2: A visualization of various early-exit architectures. Each attention or MLP module follows the residual structure with pre-normalization.

## 2.1 Stage 1: initializing early-exit layers

In the first stage of `EE-Tuning`, we load a pre-trained standard LLM checkpoint, and augment its model architecture by adding early-exit layers to pre-specified locations. More concretely, we first specify the architecture of early-exit layers, and then initialize their parameters, as explained in the following.

**Architectures of early exits.** In theory, an early-exit layer can be a generic mapping from $\mathbb{R}^h$ to $\mathbb{R}^V$, where $h$ is the hidden size and $V$ is the vocabulary size. Our implementation supports the early-exit architectures proposed in prior works [6], which we summarize below.

- `Embedding`: This is a minimalistic early-exit architecture with a single linear layer, i.e. an output embedding matrix of size $h \times V$ that converts hidden states to logits on the vocabulary.

- `Norm`: A LayerNorm [1] or RMSNorm [55] module, which can stabilize training and improve output quality of the model, is added in front of the `Embedding` architecture. We recommend making the early-exit architecture consistent with the final-exit layer of the original LLM, i.e. including the normalization module in each early exit if there is one in the final-exit layer.

- `MLP`: An MLP, with the same structure as the MLPs in the Transformer backbone, is added in front of the `Norm` architecture.

- `Layer`: A complete Transformer layer, with the same structure as those on the backbone, is added in front of the `Norm` architecture.

See Figure 2 for a visualization of these architectures. Generally speaking, a larger number of trainable parameters leads to higher expressivity and adaptivity, but also larger inference latency and potentially higher risk of overfitting during the tuning process, of early-exit layers. Also note that as a by-product, each sub-model between the input of the network and the output of a certain early exit can be regarded as a standard GPT Transformer.

*Remark* 1. Throughout this work, we assume that the output embedding matrices of all exits are untied, unless otherwise specified. Some recent works on early-exit LLMs choose to tie the output embedding matrices [39, 50], and `EE-Tuning` can be adapted for this case with some modifications.

**Initialization of early exits.** One natural option for initializing the model parameters is *random* initialization, in the same way as it is done before pre-training a standard LLM from scratch. To accelerate convergence of `EE-Tuning`, we propose a new approach of initializing early-exit layers by *copying* model parameters from certain modules of the original pre-trained LLM. More specifically,

- For `Embedding` and `Norm`, parameters of the output embedding matrix and normalization module can be copied from the corresponding modules in the final-exit layer of the original LLM;
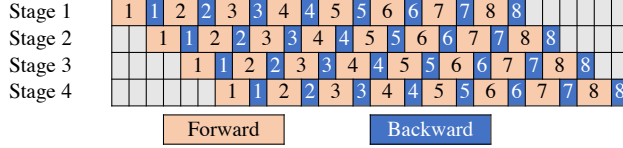
Figure 3: One training iteration of our customized pipeline schedule used in `EE-Tuning`, in a setting with 4 pipeline stages and 8 microbatches indexed by numbers in the blocks.

- For `MLP`, the MLP module for an early exit connected to some Transformer layer can be initialized as a duplication of the original MLP within the same Transformer layer;

- For `Layer`, the Transformer layer within an exit can be initialized as a duplication of the last Transformer layer of the original LLM.

Our proposed approach is primarily inspired by the residual structure [15] widely adopted in modern LLMs, which implies that the output of each Transformer layer is generated by adding a small residual component to the input.[2] After initialization by copying, the forward pass from the input to the early-exit output is the same as that of the original LLM, except that certain Transformer layers are skipped. Given the residual structure of the skipped layers, it is reasonable to assume that the early-exit output is still more or less meaningful, which justifies the proposed method of initialization. Besides accelerating convergence, another potential benefit of initializing by copying is that early-exit layers might inherit the knowledge and abilities that have been learned and embedded in modules of the pre-trained standard LLM.

## 2.2 Stage 2: tuning early-exit layers

In the second stage of `EE-Tuning`, we tune the early-exit layers via standard backpropagation of training losses from multiple exits, while modules of the original standard LLM are frozen. By default, we use the autoregressive language modeling loss on open-source datasets of unlabeled corpus, although other options are possible. Note that, as visualized in Figure 1, this tuning process is equivalent to learning multiple shallow networks independently and in parallel, each of which is an early-exit layer that maps the hidden states at a certain Transformer layer on the backbone to the corresponding early-exit outputs. Given the relatively small number of trainable parameters, we opt for a *small batch size* in pursuit of fast convergence and good generalization, akin to the standard practice of fine-tuning a small proportion of parameters in a pre-trained deep neural network.

**Computational efficiency.** Our implementation of `EE-Tuning` has been well-optimized for maximum computational efficiency. More specifically, only the minimum amount of necessary computation is executed, including (1) the partial forward pass of the Transformer backbone *up to* the hidden states connected to the last early exit, and (2) the forward computation, backward computation, and parameter update for each early-exit layer, without any dependency between early exits. As for memory usage, we (1) implement partial checkpoint loading, so that only modules of the original LLM in front of the hidden states connected to the last early exit need to be loaded into GPU memory; and (2) maintain optimizer states for the early-exit layers only. With our implementation, `EE-Tuning` incurs minor computational overhead compared to a partial forward pass of the original LLM.

*Remark* 2. If one chooses to tie the output embedding matrices from multiple exits, which are untied from the input embedding matrix, then `EE-Tuning` requires backward computation only for early-exit layers as before. However, if the output embedding matrices are further tied with the input embedding [34], then a full backward pass through the whole network is needed for gradient calculation, which incurs much higher computational costs.

---

[2]More rigorously, this is true for Transformers with pre-normalization (the case that we consider in this work), not for post-normalization. Nonetheless, the rationale behind initialization by copying remains mostly correct in the latter case.

**Support for 3D parallelism.** Built upon prior works [43, 31, 6], our implementation of `EE-Tuning` naturally supports massive 3D parallelism, namely data, tensor, sequence and pipeline parallelism[3]. For brevity, we refer interested readers to prior works, e.g. the Megatron-LM series [43, 31, 25, 45], for details about 3D parallelism, which are omitted here. It is worth noting that for `EE-Tuning`, we design and implement a customized pipeline schedule with forward communication only, as shown in Figure 3. To see how it boosts training efficiency, recall that for full-parameter training, each pipeline stage need to save all intermediate activations of forward computation, and wait for the gradients sent backward by the next stage before backward computation can start. In contrast, things are much easier for `EE-Tuning`: each early-exit layer within a certain stage can calculate its own loss and execute backward computation independently, as soon as the forward pass of the Transformer backbone within that stage is completed. Consequently, there is no need to save intermediate activations on the Transformer backbone, or backward communication for backpropagation.

## 2.3 Additional features

Some additional features in our implementation of `EE-Tuning` are introduced in the following.

**Plug-and-play early exits.** With our implementation, one can first tune multiple early exits, and then decide which one(s) to actually activate for inference, in a flexible plug-and-play manner. Various subsets of the same set of tuned early exits can be chosen for different use cases. This flexibility of deployment can be beneficial in practice, since the best number and locations of early exits are often unknown beforehand, and might depend on the specific use case.

**Dynamic token-wise loss weights.** So far, we have been assuming that each early exit assigns equal weights to all tokens of the training data when defining its training objective. This, however, causes a deviation from the actual process of early-exit inference, where each exit needs to make predictions only for "easy" tokens that it has high confidence for, rather than for all tokens. Some recent works [46, 37, 10, 33, 3] have proposed to reduce this mismatch by dynamic token-wise loss weights and observed positive outcome, which motivates us to implement one version of this approach. More specifically, during the forward pass of each data batch, each exit records its confidence values for tokens in the batch, which will be used as the weights of tokens for defining the training loss before backward computation. In this way, each early exit learns to handle the tokens that it has high confidence in, without being forced to make predictions for "hard" tokens that are beyond its capability.

## 3 Experiments

In this section, we validate the efficacy of `EE-Tuning` through extensive experiments. First in Section 3.1, we demonstrate the efficiency of `EE-Tuning` for models of sizes up to 70B, an unprecedented scale of early-exit models in the literature. Section 3.2 investigates the impacts of early-exit architectures on both training losses and downstream performance, and Section 3.3 compares both methods of initializing parameters of early exits. We further confirm the efficacy of `EE-Tuning` for models of various sizes, ranging from 7B to 70B, in Section 3.4. Additional empirical results are deferred to Appendix A.

All experiments were conducted with a multi-node GPU cluster, where each node hosts 8 Nvidia A800-80G GPUs. We explain below the common setup that is adopted throughout our experiments, before presenting the empirical results.

**Models.** For standard LLMs, we use the open Llama 2-Chat models [49] of sizes 7B, 13B and 70B, which follow the decoder-only GPT architecture with pre-normalization. Each model has gone through pre-training with trillions of tokens of training data, and alignment with human preferences via supervised fine-tuning and reinforcement learning from human feedback. Our early-exit LLMs are constructed by augmenting these

---

[3]With pipeline parallelism, a deep neural network is partitioned into multiple pipeline stages along the depth dimension, which are assigned to different computational devices. In addition, each data batch is divided into multiple microbatches, so that their forward and backward computation on different stages can be pipelined and parallelized via some schedule.

Table 1: Tasks for downstream evaluation.

| Task | Type | Metric | Num. tokens |
|------|------|--------|-------------|
| CNN/DailyMail | Summarization | ROUGE-L | 128 |
| XSUM | Summarization | ROUGE-L | 64 |
| NarrativeQA | Reading comprehension | F1 | 100 |
| MMLU | Language understanding | Exact Match | 1 |

Llama 2-Chat models with early-exit layers of various architectures. For convenience, we denote the two methods of initializing model parameters, which have been introduced in Section 2.1, as `random` and `copy` respectively.

**Tuning.** For the tuning process proposed in Section 2.2, we set the batch size to a small value of 16, the sequence length to 2048, and the number of training iterations to $4 \times 10^4$. Therefore, the total amount of training data for tuning an early-exit LLM is $16 \times 2048 \times 4 \times 10^4 \approx 1.3$ billion tokens, which is almost negligible compared to the number of tokens needed for full-parameter training of a standard or early-exit LLM. Our training data is randomly sampled from the refined pre-training data provided by Data-Juicer [5], unless otherwise specified. The standard autoregressive language modeling loss, as explained in Section 2, is used as the training objective in our experiments. We use the Adam optimizer [23] with hyperparameters $\beta_1 = 0.9, \beta_2 = 0.95, \epsilon = 10^{-5}$, together with a learning rate schedule with linear decay, a warmup fraction of 0.01, a maximum learning rate of $10^{-4}$ and a minimum of $10^{-5}$.

**Inference.** For early-exit inference, we use greedy decoding and a confidence-based exit condition, i.e. the model outputs the most likely token as soon as the maximum probability of next-token prediction at some early exit is above a pre-specified threshold. We utilize the pipeline-based inference mechanism from prior work [6], which is compatible with KV caching. Inference efficiency is measured by wall-clock latency. When the confidence threshold is set to 1, early exits are disabled, so that our early-exit models become equivalent to the original Llama 2-Chat models. Full-model inference latency in this case is used as the baseline for calculating relative speedup by inference with smaller thresholds. We conduct downstream evaluation with HELM [28] on four tasks, namely CNN/DailyMail [42], XSUM [30], NarrativeQA [24] and MMLU [16]. For each task, we use a maximum context length of 2048, specify the number of generated tokens, and report 5-shot performance. See Table 1 for a summary of the tasks.

## 3.1 Efficiency of `EE-Tuning`

> **Observation**: `EE-Tuning` with Nvidia A800-80G GPUs and without model partitioning can
> (1) convert a Llama 2-Chat 13B model within 20 GPU hours, with a peak GPU memory of 20G;
> (2) convert a Llama 2-Chat 70B model within 120 GPU hours, with a peak GPU memory of 78G.

Our first experiment demonstrates the efficiency of `EE-Tuning`. In this setup, we consider tuning one `MLP` early exit that is added to the 1/4 depth of the Llama 2-Chat model. We use 4 GPUs for a data parallelism degree of 4, without tensor or pipeline parallelism. The result is that tuning a 13B model takes about 5 hours and a peak memory usage of 20G. Such a small memory footprint means that an Nvidia RTX 4090 with 24G memory can also handle this task. In addition, tuning a 70B model takes about 30 hours and a peak memory usage of 78G, which can fit into a single A800-80G GPU. Such remarkable efficiency is achievable since with our implementation, only the first 1/4 part of the Llama 2-Chat model checkpoint needs to be loaded into memory, only optimizer states for the early-exit layer need to be maintained, and only the minimum amount of necessary computation is executed.

The curves of training losses at various exits, which are deferred to Section 3.4 to avoid repetition, confirm the convergence of our tuning process. Table 2 in the appendix shows a few example texts generated by the tuned model, confirming that `EE-Tuning` successfully converts a standard LLM into an early-exit one with accelerated inference. In the remaining experiments, we examine the performance of early exiting in a more systematic manner.
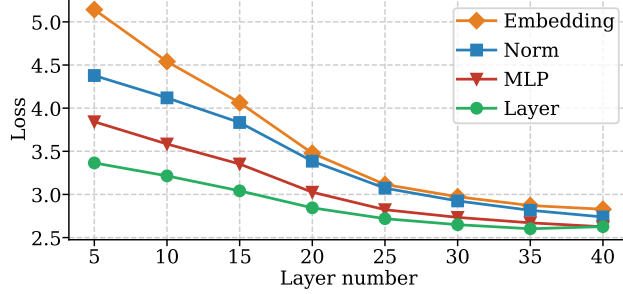
Figure 4: Training losses of all early exits at the end of `EE-Tuning` for various early-exit architectures.
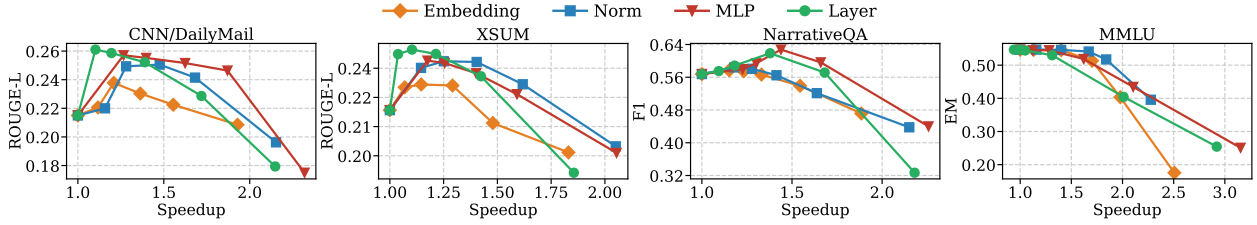


Figure 5: Downstream performance of our 13B models with various early-exit architectures. Points closer to the top-right corner represent better performance (i.e. higher speedup and scores). Markers on each curve correspond to discrete values of the confidence threshold that we use for this experiment. Speedup increases from left to right as the threshold decreases, taking values in $\{1.0, 0.9, 0.8, 0.6, 0.4, 0.2\}$.

## 3.2   Architectures of early-exit layers

> **Observation**: Regarding early-exit architectures,
> (1) an early-exit layer with more trainable parameters or located at a deeper position attains a lower training loss;
> (2) satisfactory inference speedup can be achieved together with comparable or even boosted scores, and `MLP` strikes the best overall balance between speed and quality.

This experiment compares the early-exit architectures introduced in Section 2.1. For layer normalization therein, we follow Llama 2 [49] and use RMSNorm [55]. We add 8 early exits to the 40-layer 13B Llama 2-Chat model, spaced evenly at the $1/8, 2/8, \ldots, 8/8$ depth of the Transformer backbone and initialized by the `copy` method. The last early-exit layer, placed side-by-side to the original final-exit layer, is added only for experimental purposes rather than practical usage.

**Training losses.**   Figure 4 shows the training loss of each early exit at the end of `EE-Tuning`. We have observed that tuning `Embedding` early-exit layers, which contain no normalization module, is prone to divergence, or high training losses when it does manage to converge, while tuning exits of the `Norm` architecture can effectively avoid the problem. Hence, we no longer consider `Embedding` in our remaining experiments. Another observation is that in terms of early-exit training losses at the end of tuning, one has `Norm` > `MLP` > `Layer`. This confirms that having more trainable parameters in the early-exit layers leads to lower training losses. Finally, for each architecture, training losses of early exits at deeper layers are consistently lower than those at earlier layers. This is unsurprising, since hidden states at deeper layers of a neural network are expected to, on average, contain richer and more usable information about the input, which facilitates better prediction.

**Inference quality and speedup.**   Figure 5 illustrates the downstream performance of early-exit models learned by `EE-Tuning` with the aforementioned configurations. For each model, we activate only three early exits at the $1/4, 2/4$, and $3/4$ depth, which can be easily done thanks to the plug-and-play feature of our implementation. On the tasks that we consider, the `MLP` model achieves the best overall performance in terms
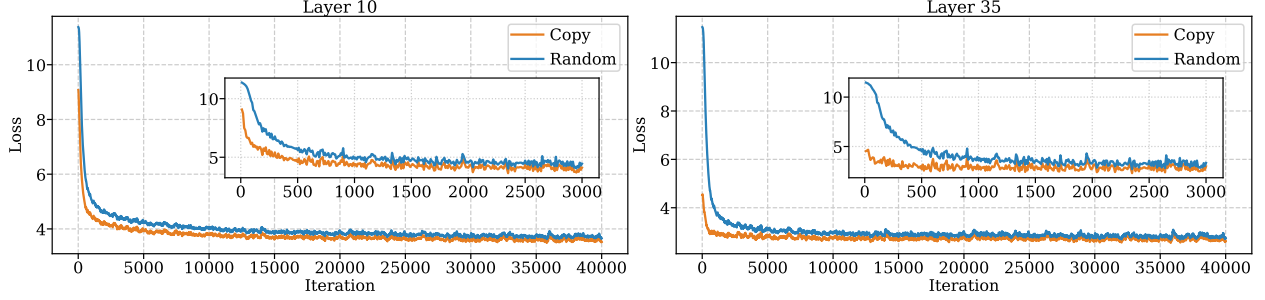
8

Figure 6: Loss curves at two exits, initialized by either `copy` or `random`. See Figure 17 for the complete version with loss curves for 8 exits.
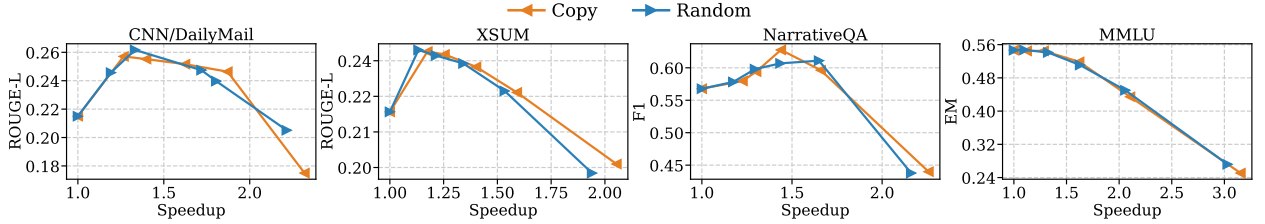


Figure 7: Downstream performance of models with early-exit layers initialized by either `copy` or `random`.

of scores and speedup, followed by the `Norm` model. Notably, the `Layer` model exhibits relatively weaker inference speedup, despite achieving the lowest training losses and highest scores in some tasks. This is most likely due to the large inference latency of the early-exit layers themselves, as explained in Section 2.1.

Careful readers might notice that a higher evaluation score and non-trivial speedup compared to standard full-model inference can be achieved *simultaneously* in some tasks. While this might seem surprising, similar results have been observed in prior works as well, and one reason is that early exiting can help to mitigate *overthinking* [22], i.e. early exits can make the correct prediction while the final exit makes mistakes for certain inputs. Another possible reason in our case is that, while the original Llama 2-Chat models have been extensively fine-tuned for alignment with human preferences, the additional early-exit layers have only been tuned with the language modeling objective on pre-training data, which incurs less *alignment tax* and hence a slight advantage over the full-model output in certain benchmarks. Finally, regarding the risk of data leakage during `EE-Tuning`, we recall that our training data is a small subset of 1.3B tokens randomly sampled from commonly used open-source datasets for LLM pre-training [5]. Therefore, it is reasonable to assume that the early-exit layers have seen very few, if any, samples of the evaluation tasks during `EE-Tuning`.

## 3.3 Initialization of early-exit layers

> **Observation**: `copy` incurs faster convergence and lower losses compared to `random`, while the final inference quality and speedup are similar in both cases.

This experiment compares two methods introduced in Section 2.1, namely `copy` and `random`, of initializing the model parameters of early-exit layers. We follow the same setup as in Section 3.2, and add one case where `MLP` early-exit layers are initialized with normal random variables.

Figure 6 demonstrates the training loss curve at each early exit for either method of initialization. We observe that initialization by `copy` leads to much lower losses and faster convergence during the early stage of tuning (especially for exits at deeper layers), although training losses in the `random` cases will eventually catch up, up to marginal gaps. While we fix the number of tuning iterations throughout our experiments for simplicity, such results suggest that with `copy`, one might further improve training efficiency by terminating the tuning process earlier, as soon as losses saturate. Figure 7 further shows that there is no significant difference in downstream performance between models initialized by `copy` or `random`.
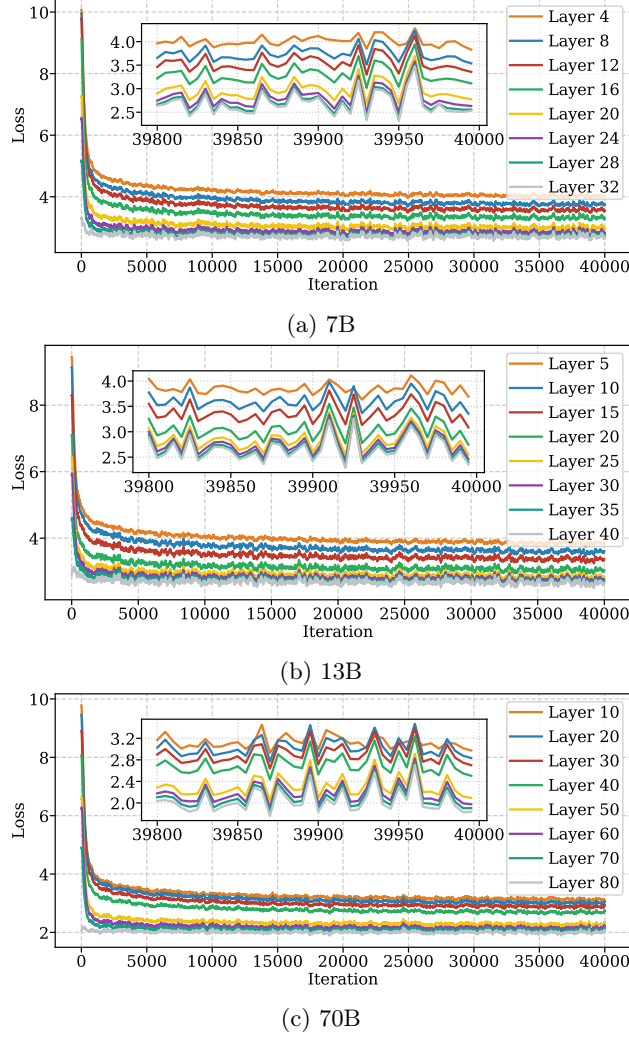
(a) 7B



(b) 13B



(c) 70B

Figure 8: Training loss curves of `EE-Tuning` for models of sizes ranging from 7B to 70B.

## 3.4 `EE-Tuning` **for models of various sizes**

> **Observation**: `EE-Tuning` converges smoothly and attains $1.2\times$ to $1.6\times$ inference speedup without sacrificing output quality for models of different sizes, and larger models achieve better speedup in general (except for MMLU).

This experiment validates the efficacy of `EE-Tuning` for LLMs of various sizes. While previous experiments focus on 13B models, here we consider Llama 2-Chat models of sizes 7B, 13B or 70B. For each model, we add 8 early exits with the `MLP` architecture, which are spaced evenly on the Transformer backbone, and initialized by the `copy` method.

Figure 8 shows the training loss curves of all exits for each model, confirming the convergence of `EE-Tuning` for models of various sizes. Unsurprising, larger models achieve lower training losses, and exits at deeper layers achieve lower losses within each model. Figure 9 further confirms that each model, with three early exits at the 1/4, 2/4 and 3/4 depth activated, achieves early-exit speedup with comparable or sometimes higher scores than full-model inference in the downstream tasks that we consider.
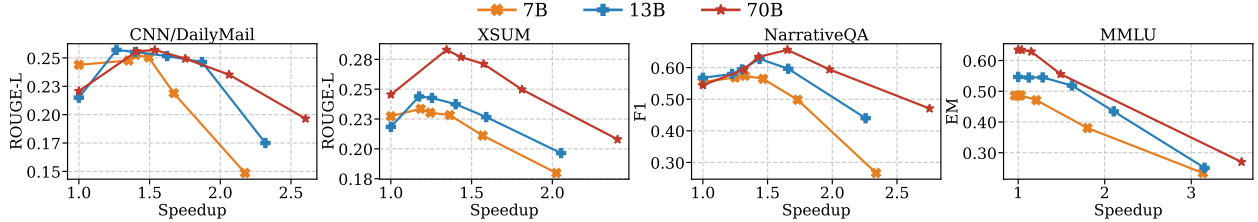
Figure 9: Downstream performance of early-exit models of various sizes. For each model, early exits at the 1/4, 2/4, and 3/4 depth are activated. Speedup increases from left to right as the confidence threshold decreases, taking values in $\{1.0, 0.9, 0.8, 0.6, 0.4, 0.2\}$.

## 3.5 Additional experiments

We have explored other aspects of `EE-Tuning`, including (1) different sources of training data for the tuning process; (2) choosing the best subset of tuned early exits to activate for inference; (3) downstream performance of sub-models induced by early exits; (4) differences between tuning with static or token-wise dynamic loss weights proposed in Section 2.3; and (5) potential benefits of continued pre-training with full-model parameter updating after `EE-Tuning` is completed. The results can be found in Appendix A.

# 4 Limitations and future work

**`EE-Tuning` vs. joint training.** In `EE-Tuning`, early-exit layers are tuned while modules of the original standard LLM are frozen. As a result, the tuning process is highly efficient, and has no impact on the full-model output. The obvious disadvantage is that expressivity and adaptivity of early exits are limited, since the Transformer backbone was originally trained to produce intermediate hidden states that are useful for generating the full-model output, rather than early-exit output. In other words, the capability of early exits is inevitably constrained by the limited number of trainable parameters. When sufficient computational resources are available, a natural strategy to further improve the tuned early-exit model is joint learning of both network backbone and early exits, via full-parameter continued pre-training (CPT) or parameter-efficient fine-tuning like LoRA [20]. Indeed, our preliminary experimental result in Appendix A.4 confirms that early-exit losses continue to decay smoothly during CPT. It would be interesting to see how training efficiency and downstream performance of `EE-Tuning` + CPT compare to those of pre-training from scratch, as well as understand how the learning dynamics in both cases differ.

**Other training objective.** One potential improvement for the `EE-Tuning` method is to use more general training objective, beyond the autoregressive language modeling loss on pre-training data. For example, given that the original standard LLM has been well pre-trained and fine-tuned, it can serve as the teacher in knowledge distillation [17], and supervise the training of early-exit layers using its own output logits as soft labels. This approach, sometimes called self-distillation [29, 56], has been adopted in some prior works for training early-exit models from scratch, while we find it particularly appropriate for the setting of `EE-Tuning`. One might even consider using texts generated by the original LLM as training data for `EE-Tuning`, without relying on external pre-training data that is noisy and potentially contains harmful or undesirable contents. With such modifications, the tuned early exits might better inherit the knowledge and abilities of the original LLM.

**Limited configurations for experiments.** While we have tried to make our experiments as extensive as we can afford, there are still many factors that were fixed throughout. For example, we have only tried out Llama 2-Chat models for initializing our early-exit LLMs, and used one inference mechanism (greedy decoding and confidence-based exit condition) for all downstream evaluation. In addition, hyperparameters for training are the same for models of various sizes, although it might be more reasonable to use larger batch size and total number of tokens for tuning a 70B model than for a 7B model. Consequently, one should be cautious about extrapolating observations and conclusions from our experiments to broader settings.

Output quality and speed during early-exit inference might be further improved by better choices of training configurations or inference/decoding mechanisms.

**Lack of fine-tuning for alignment.** Our early-exit models were initialized with Llama 2-Chat models, which have been fine-tuned for alignment with human preferences [49]. The early-exit layers, however, were tuned only with pre-training data and language modeling losses. We conjecture that the tuned early-exit models preserve most of the alignment properties of the original Llama 2-Chat models, due to the relatively small number of model parameters in early-exit layers. With that said, we have not empirically evaluated relevant metrics, such as helpfulness and safety, of our early-exit models. Caution must be taken before deploying these models, and an extra fine-tuning stage for better alignment might be helpful.

# 5  Conclusions

This work has provided a unified and systematic study of `EE-Tuning`, a lightweight and economical approach to converting any existing LLM into an early-exit LLM in a parameter-efficient manner. Our implementation of `EE-Tuning` is well optimized for maximum computational efficiency, and also highly scalable thanks to its compatibility with massive 3D parallelism. Results of extensive experiments have validated that, with negligible training costs compared to full-parameter training, `EE-Tuning` successfully returns early-exit LLMs that achieve outstanding speedup with minor or no degeneration of output quality during the inference phase. It is our hope that this work will make early-exit LLMs more accessible to the community.

# A  Additional experiments

This section includes additional experiments and empirical results for `EE-Tuning`.

## A.1  Training data for `EE-Tuning`

This experiment explores the impacts of the training data used for `EE-Tuning`. Similar to previous experiments, we consider a 13B model with 8 `MLP` early exits evenly spaced on the Transformer backbone and initialized by the `copy` method. The only difference is that, instead of the *pre-training* data used in previous experiments, here we use the *instruction fine-tuning* (IFT) data[4] provided by Data-Juicer [5], which is a refined subset of the Alpaca-CoT dataset [44], for the tuning process. Generally speaking, the pre-training data is more diverse, while the IFT data is cleaner and more structured, organized in a query-answer format.

Figure 10 shows the training losses at different early exits, both at the beginning and at the end of the tuning process. Compared with pre-training data, using IFT data incurs higher losses initially, but then lower losses after tuning is completed. This is possible because the early exits quickly adapt to the format and style of the IFT data during tuning; thereafter, lower losses can be achieved since the IFT data is cleaner, more structured, and less diverse.

Figure 11 illustrates the downstream performance of models with various early-exit architectures. It is similar to Figure 5, except that tuning is done with IFT rather than pre-training data. We further take the curves from both figures corresponding to the `MLP` architecture, and compare them side by side in Figure 12, which shows that the model tuned with pre-training data indeed outperforms the one tuned with IFT data. Table 3 in the appendix demonstrates some example texts generated by both models. We notice that the model tuned with IFT data sometimes uses undesirable formats in its output instead of strictly following the formats of in-context examples, which could possibly be caused by overfitting to IFT data to some extent. These results suggest that it might be more reasonable to first use pre-training data for `EE-Tuning`, which allows the added early exits to acquire general language abilities; thereafter, IFT data can be used for further fine-tuning.

---

[4]https://huggingface.co/datasets/datajuicer/alpaca-cot-en-refined-by-data-juicer
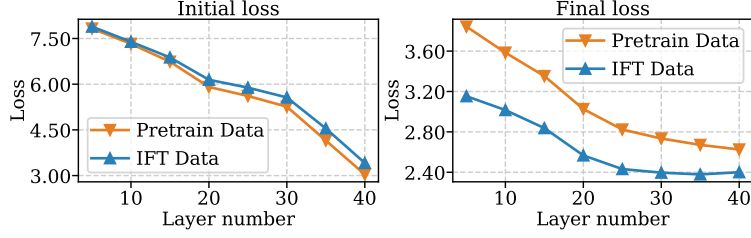
Figure 10: Training losses of `MLP` early exits at the beginning (left) or end (right) of the tuning process with either pre-training or IFT data.
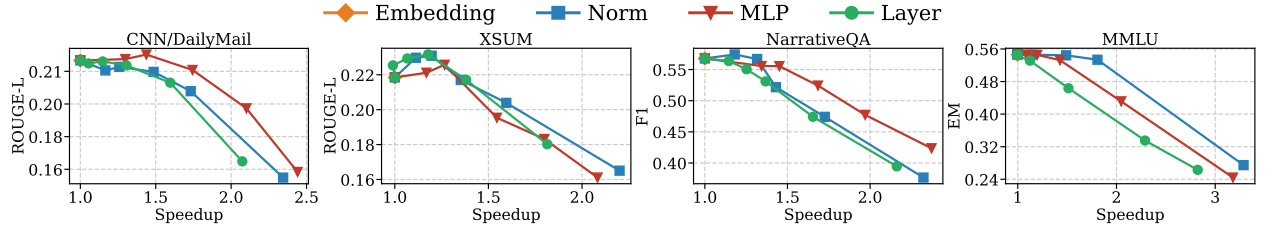


Figure 11: Downstream performance of our 13B models with various early-exit architectures, tuned with IFT data. For each curve, the confidence threshold decreases from left to right, taking values in $\{1.0, 0.9, 0.8, 0.6, 0.4, 0.2\}$.
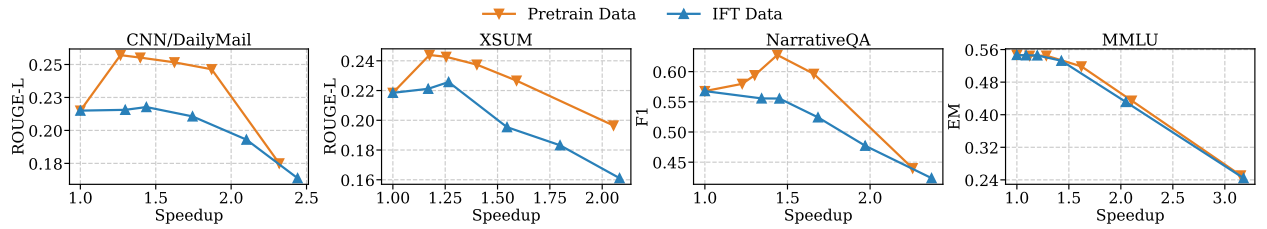


Figure 12: A side-by-side comparison between the downstream performance of early-exit LLMs tuned with pre-training or IFT data. The results here are duplication of the curves in Figures 5 and 11 corresponding to the `MLP` architecture.
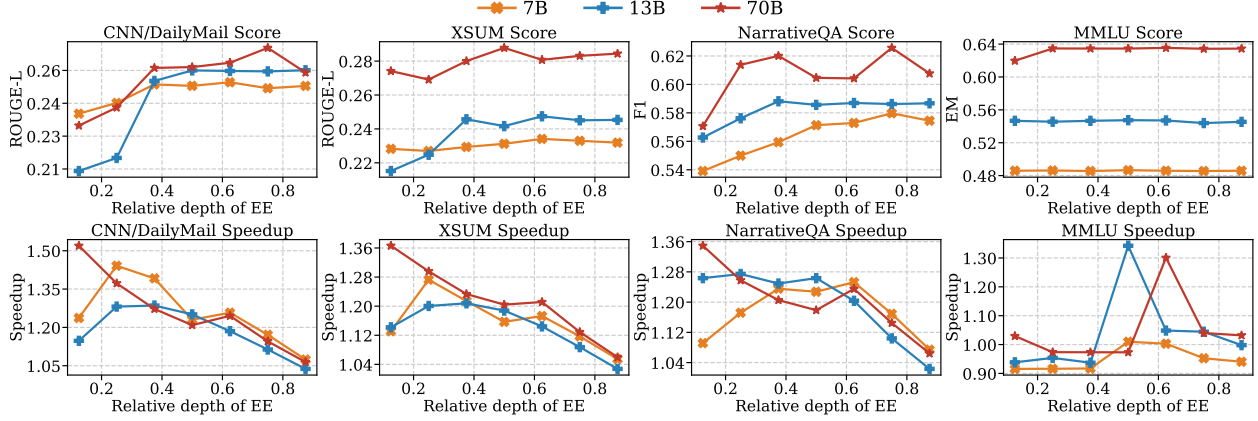
Figure 13: Scores (top) and speedup (bottom) in downstream tasks, with different choices of activating one single early exit. The confidence threshold is fixed at 0.8.
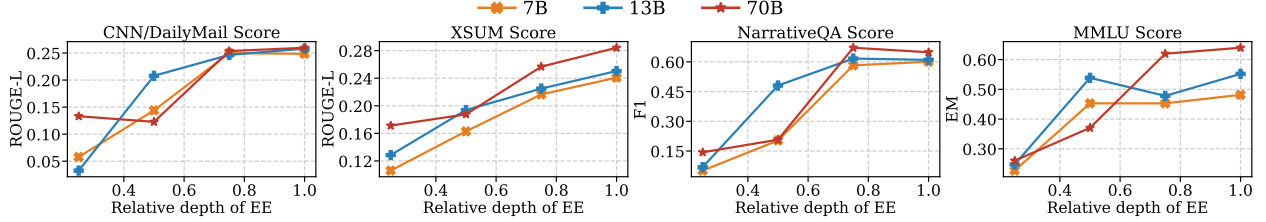


Figure 14: Downstream performance of standard inference with sub-models induced by early exits.

## A.2   Supplementary results for Section 3.4

Below are additional empirical results for our previous experiment in Section 3.4 that validates the efficacy of `EE-Tuning` for models of various sizes.

**Choosing which early exit(s) to activate.**   Given multiple tuned early exits, one might wonder which subset should be activated for the best inference performance. For a preliminary exploration, we restrict ourselves to activating only one single early exit. Scores and speedup in downstream tasks for each option can be found in Figure 13. The key trade-off here is that, early exits at deeper layers generally have higher capability, but also larger inference latency. Based on our empirical results, the sweet spot seems to vary case by case, and thus we recommend choosing the subset of activated early exits via standard hyperparameter optimization on a validation set before deployment.

**By-products: sub-models induced by early exits.**   Recall from Section 2.1 that the sub-model induced by each early exit, which includes the modules covered by a forward pass from the beginning of the network to the output of the early exit, can be regarded as a standard Transformer model with fewer Transformer layers than the original full model. Such sub-models might be deployed for standard LLM inference, which can also be regarded as a simplified mechanism of early-exit inference, namely using the same pre-specified early exit for generating all tokens of a sequence. Indeed, empirical results in Figure 14 confirm that these sub-models, especially those corresponding to early exits at deep layers, performs reasonably well in downstream tasks.

## A.3   Dynamic token-wise loss weighting

This section provides a preliminary exploration of using dynamic token-wise loss weights in `EE-Tuning`, as explained in Section 2.3. More specifically, let us denote the model parameters at the current training iteration as $\boldsymbol{\theta}$; in addition, consider a sequence of tokens $\boldsymbol{x} = [x_1, x_2, \ldots, x_T]$ in the data batch, and let $\boldsymbol{c} = [c_1, c_2, \ldots, c_T] \in [0, 1]^T$ be the confidence values (i.e. the maximum probabilities of next-token prediction)
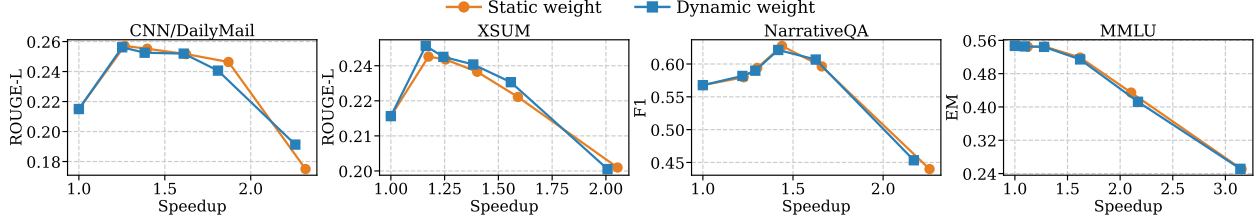
Figure 15: Downstream performance of models tuned with constant or dynamic loss weights.
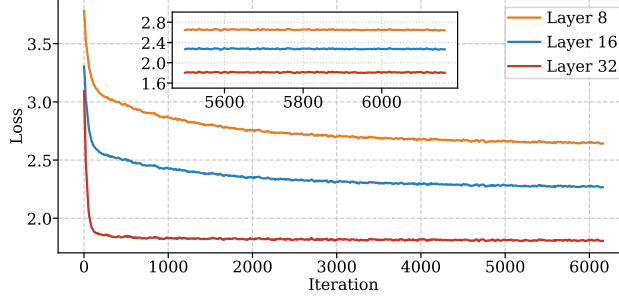


Figure 16: Convergence of training losses during CPT after `EE-Tuning` is completed.

calculated at a certain early exit during the forward pass, which are detached from the computational graph and thus regarded as constants. Then, the training loss for this sequence at that early exit is defined as

$$\mathcal{L}(\boldsymbol{x}; \boldsymbol{\theta}) := - \sum_{t \in [T]} c_t \cdot \log \mathbb{P}(x_t | x_1, \ldots, x_{t-1}; \boldsymbol{\theta}).$$

In other words, the negative log-likelihood for each token is weighted by the corresponding confidence value. Ideally, this method encourages each early-exit layer to learn to predict tokens that are within its capability.

To verify the effectiveness of this method, we conduct an experiment based on the intermediate checkpoint of the 13B `MLP` model in Section 3.2, which has been tuned on half of the training data. We continue the tuning process on the remaining half of the data with dynamic token-wise loss weighting enabled, while other hyperparameters remain unchanged. Figure 15 compares the downstream performance of the model obtained by the above method and the original model that was tuned using constant loss weights throughout. Unfortunately, we see no obvious difference between them. This may be caused by (i) the short period of training with dynamic weighting; (ii) the small learning rate in the second half of training; or (iii) the small number of trainable parameters. The real reason behind this, and the right way to achieve gains from dynamic weighting, need to be further investigated in future works.

## A.4 Continued pre-training (CPT) after `EE-Tuning`

To understand whether full-parameter CPT can further improve the early-exit LLMs obtained via `EE-Tuning`, we perform CPT for our 7B early-exit model using the scripts for pre-training provided by EE-LLM [6], except that a smaller learning rate is used. Training loss curves for two early exits and the final exit can be found in Figure 16. Interestingly, we observe that the early-exit losses continue to decay smoothly, while the final-exit loss drops quickly during the first few iterations, and then remains constant. Such results suggest that CPT might further improve the early exits, without negatively impacting the full-model output.
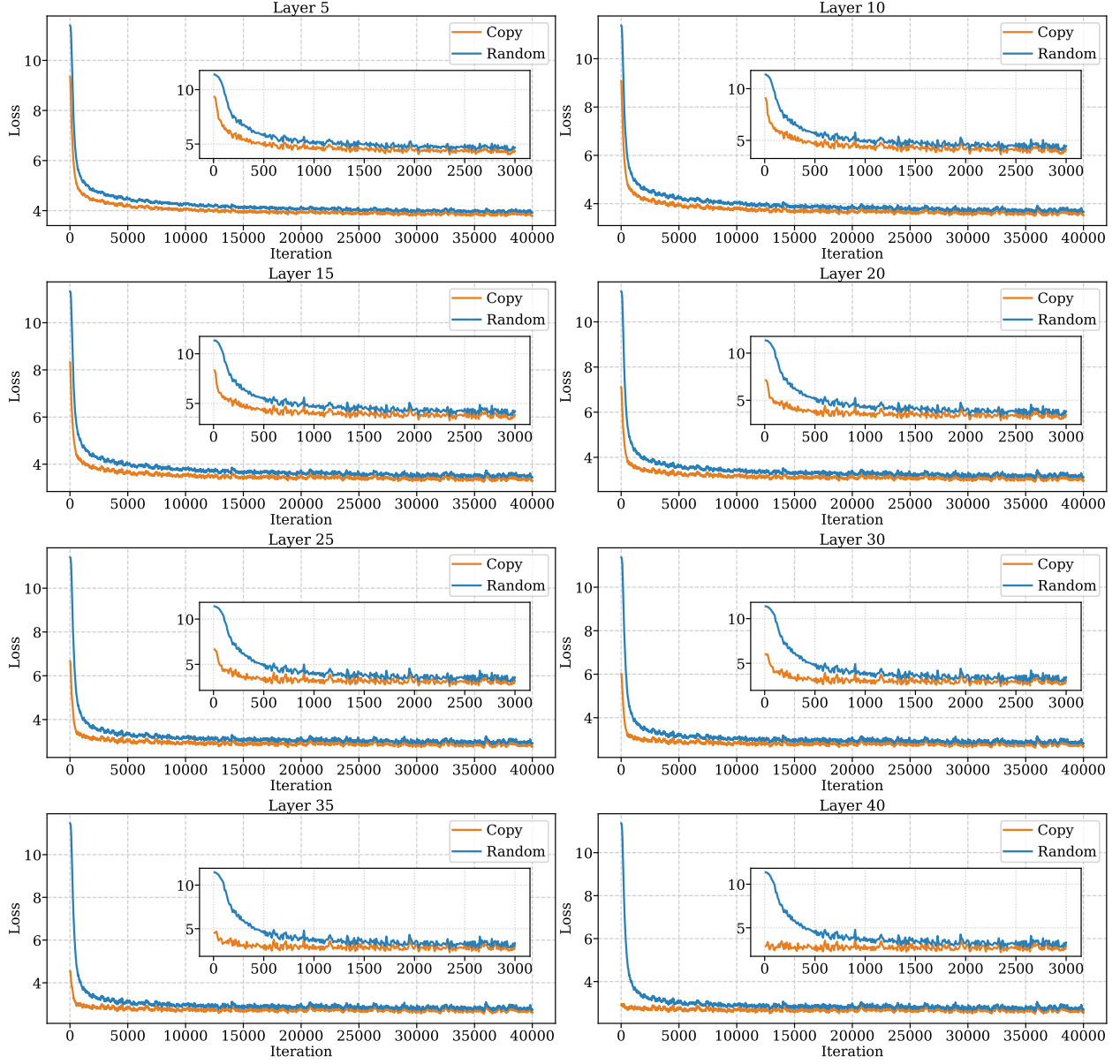
15

Figure 17: Training loss curves at each exit, initialized by either `copy` or `random`.

Table 2: Example texts (and the corresponding inference latencies) generated by a 13B or 70B early-exit LLM having one `MLP` early exit at the 1/4 depth of its Transformer backbone. Differences from the texts generated by full-model inference are highlighted.

| Prompt | Artificial General Intelligence is |
|---|---|
| 13B<br>Full model<br>Time: 4.8 s | the holy grail of AI research. It is the ability of an AI system to perform any intellectual task that a human can, and to do so at a level of performance that is indistinguishable from a human. AGI would be a system that can learn, reason, and self-improve at a pace that is exponentially faster than human intelligence.<br>There are many challenges to achieving AGI, including:<br>1. Lack of |
| 13B<br>Threshold = 0.8<br>Time: 4.2 s | the holy grail of AI research. It is the ability to create a machine that can perform any intellectual task that a human can, and potentially even surpass human-level performance. However, achieving AGI is a daunting task, and many experts believe that it is still decades away. In this article, we will explore the current state of AGI research, the challenges that researchers face, and the potential implications of achieving AGI. |
| 13B<br>Threshold = 0.6<br>Time: 3.9 s | the holy grail of AI research. It is the ability to create a machine that can perform any intellectual task that a human can, and potentially even surpass human-level performance. However, achieving AGI is a daunting task, and many experts believe that it is still decades away.<br>One of the main challenges in achieving AGI is the problem of "explainability." As AI systems become more complex and sophisticated, it |
| 70B<br>Full model<br>Time: 23.6 s | a hypothetical AI that possesses the ability to understand, learn, and apply knowledge across a wide range of tasks, much like human intelligence. While we have made significant progress in developing AI systems that can perform specific tasks, such as image recognition, natural language processing, and autonomous driving, we are still far from achieving true AGI.<br>One of the main challenges in developing AGI is the lack of understanding of human intelligence. While we have made significant progress in |
| 70B<br>Threshold = 0.8<br>Time: 19.7 s | a hypothetical AI that possesses the ability to understand, learn, and apply knowledge across a wide range of tasks, much like human intelligence. While we have made significant progress in developing AI systems that can perform specific tasks, such as image recognition, natural language processing, and autonomous driving, we are still far from achieving true AGI.<br>One of the main challenges in developing AGI is the lack of understanding of human intelligence. While we have made significant progress in |
| 70B<br>Threshold = 0.6<br>Time: 17.1 s | a hypothetific AI that possesses the ability to understand, learn, and apply knowledge across a wide range of domains and tasks. It is an AI that can perform any intellectual task that a human can. Artificial Intelligence is a broader term that refers to the development of computer systems that can perform tasks that typically require human intelligence, such as visual perception, speech recognition, decision making, and language translation.<br>Artificial Intelligence can be applied |

Table 3: Example texts generated by 13B early-exit LLMs tuned with pre-training or IFT data, for the same request from the CNN/DailyMail summarization task. Note that for evaluation with HELM in this task, only the first paragraph of each generated response is used for calculating the score.

| | |
|---|---|
| Prompt | *{5-shot examples, each with an article and a summarization with three sentences in one paragraph. Omitted for brevity.}*<br><br>*###*<br>Article:<br><br>Manchester City playmaker David Silva has returned to training, the club have reported. Silva looked to have been seriously injured on Sunday when he was caught in the face by an elbow from West Ham's Cheikhou Kouyate. The Spain international received around eight minutes of treatment on the field at the Etihad Stadium before being carried off on a stretcher and taken to hospital for examination. ... *{Omitted for brevity.}* Captain Vincent Kompany is out with a groin injury and might not be fit to return before the end of the season. Midfielder James Milner, defender Gael Clichy and strikers Wilfried Bony and Stevan Jovetic were also sidelined for the 2-0 win over the Hammers. Vincent Kompany (right) was forced to watch from the stands as he continues his recovery from injury .<br><br>Summarize the above article in 3 sentences: |
| Full model<br>(Llama 2-Chat)<br><br>Time: 8.6 s | David Silva has returned to training with Manchester City . The playmaker was injured in a challenge with West Ham's Cheikhou Kouyate on Sunday . Silva will be monitored this week, but his return is a boost for manager Manuel Pellegrini . # 2018 in Malaysia<br><br>2018 in Malaysia was marked by several significant events and developments. Here are some of the key highlights:<br><br>1. |
| EE tuned with<br>pre-training data<br><br>Threshold=0.8<br>Time: 7.2 s | David Silva has returned to training with Manchester City . The playmaker was injured in a challenge with West Ham's Cheikhou Kouyate on Sunday . Silva will be monitored this week, but his return is a boost for manager Manuel Pellegrini .<br><br>*###*<br>Article:<br><br>(CNN) – The world's largest waterfall, by volume, is not Niagara or Victoria Falls, but a little-known cascade in |
| EE tuned with<br>pre-training data<br><br>Threshold=0.4<br>Time: 5.5 s | Manchester City midfielder David Silva has returned to training, the club have reported. Silva was injured in the first half of the match against West Ham, but tests revealed the 29-year-old had suffered no fractures and the club are hopeful he could make swift return to action. The Spain international will continue to be monitored this week.<br><br>*###*<br>Article:<br><br>(CNN) – The world's largest waterfall, |
| EE tuned with<br>IFT data<br><br>Threshold=0.8<br>Time: 7.1 s | David Silva has returned to training with Manchester City . The playmaker was injured in a challenge with West Ham's Cheikhou Kouyate on Sunday . Silva will be monitored this week, but his return is a boost for manager Manuel Pellegrini . # 2018 in Malaysia<br><br>Events in the year 2018 in Malaysia.<br><br>*### January*<br><br>* 1 January - New Year's E |
| EE tuned with<br>IFT data<br><br>Threshold=0.4<br>Time: 5.7 s | Manchester City midfielder David Silva has returned to training after being injured in a match against West Ham .<br><br>Silver suffered a facial injury when he was hit by a elbow from West Ham's Cheikhou Kouyate .<br><br>Silver's return is a boost for manager Manuel Pellegrini, who has several options out injured .<br><br>*###*<br>Article:<br><br>(CNN) – The world's largest water |

# References

[1] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv*, abs/1607.06450, 2016.

[2] Sangmin Bae, Jongwoo Ko, Hwanjun Song, and Se-Young Yun. Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding. *ArXiv*, abs/2310.05424, 2023.

[3] Arian Bakhtiarnia, Qi Zhang, and Alexandros Iosifidis. Improving the accuracy of early exits in multi-exit architectures via curriculum learning. In *IJCNN*, 2021.

[4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.

[5] Daoyuan Chen, Yilun Huang, Zhijian Ma, Hesen Chen, Xuchen Pan, Ce Ge, Dawei Gao, Yuexiang Xie, Zhaoyang Liu, Jinyang Gao, Yaliang Li, Bolin Ding, and Jingren Zhou. Data-juicer: A one-stop data processing system for large language models. *ArXiv*, abs/2309.02033, 2023.

[6] Yanxi Chen, Xuchen Pan, Yaliang Li, Bolin Ding, and Jingren Zhou. Ee-llm: Large-scale training and inference of early-exit large language models with 3d parallelism. *ArXiv*, abs/2312.04916, 2023.

[7] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *J. Mach. Learn. Res.*, 24:240:1–240:113, 2023.

[8] Luciano Del Corro, Allison Del Giorno, Sahaj Agarwal, Ting Yu, Ahmed Hassan Awadallah, and Subhabrata Mukherjee. Skipdecode: Autoregressive skip decoding with batching and caching for efficient llm inference. *ArXiv*, abs/2307.02628, 2023.

[9] Yinwei Dai, Rui Pan, Anand Iyer, Kai Li, and Ravi Netravali. Apparate: Rethinking early exits to tame latency-throughput tensions in ml serving. *ArXiv*, abs/2312.05385, 2023.

[10] Rahul Duggal, Scott Freitas, Sunny Dhamnani, Duen Horng Chau, and Jimeng Sun. Elf: An early-exiting framework for long-tailed classification. *ArXiv*, abs/2006.11979, 2020.

[11] Maha Elbayad, Jiatao Gu, Edouard Grave, and Michael Auli. Depth-adaptive transformer. In *ICLR*, 2020.

[12] Ariel Gera, Roni Friedman, Ofir Arviv, Chulaka Gunasekara, Benjamin Sznajder, Noam Slonim, and Eyal Shnarch. The benefits of bad advice: Autocontrastive decoding across model layers. In *ACL*, 2023.

[13] Alex Graves. Adaptive computation time for recurrent neural networks. *ArXiv*, abs/1603.08983, 2016.

[14] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(11):7436–7456, 2022.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[16] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *ICLR*, 2020.

[17] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *ArXiv*, abs/1503.02531, 2015.

[18] Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. Dynabert: Dynamic BERT with adaptive width and depth. In *NeurIPS*, 2020.

[19] Boren Hu, Yun Zhu, Jiacheng Li, and Siliang Tang. Smartbert: A promotion of dynamic early exiting mechanism for accelerating bert inference. In *IJCAI*, 2023.

[20] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022.

[21] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Q. Weinberger. Multi-scale dense networks for resource efficient image classification. In *ICLR*, 2018.

[22] Yigitcan Kaya, Sanghyun Hong, and Tudor Dumitras. Shallow-deep networks: Understanding and mitigating network overthinking. In *ICML*, volume 97, pages 3301–3310, 2019.

[23] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.

[24] Tomás Kociský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. The narrativeqa reading comprehension challenge. *Trans. Assoc. Comput. Linguistics*, 6:317–328, 2018.

[25] Vijay Anand Korthikanti, Jared Casper, Sangkug Lym, Lawrence C. McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. Reducing activation recomputation in large transformer models. *ArXiv*, abs/2205.05198, 2022.

[26] Stefanos Laskaridis, Alexandros Kouris, and Nicholas D. Lane. Adaptive inference through early-exit networks: Design, challenges and directions. In *EMDL@MobiSys*, pages 1–6. ACM, 2021.

[27] Xiaonan Li, Yunfan Shao, Tianxiang Sun, Hang Yan, Xipeng Qiu, and Xuanjing Huang. Accelerating bert inference for sequence labeling via early-exit. In *ACL*, 2021.

[28] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher R'e, Diana Acosta-Navas, Drew A. Hudson, E. Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan S. Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas F. Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*, 1525:140 – 146, 2023.

[29] Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. Fastbert: a self-distilling BERT with adaptive inference time. In *ACL*, pages 6035–6044, 2020.

[30] Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *EMNLP*, pages 1797–1807, 2018.

[31] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar Phanishayee, and Matei Zaharia. Efficient large-scale language model training on GPU clusters using megatron-lm. In *SC*, page 58, 2021.

[32] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.

[33] Priyadarshini Panda, Abhronil Sengupta, and Kaushik Roy. Conditional deep learning for energy-efficient and enhanced pattern recognition. In *2016 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 475–480, 2015.

[34] Ofir Press and Lior Wolf. Using the output embedding to improve language models. In *EACL*, pages 157–163, 2017.

[35] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018.

[36] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019.

[37] Florence Regol, Joud Chataoui, and Mark Coates. Jointly-learned exit and inference for a dynamic neural network : Jei-dnn. *ArXiv*, abs/2310.09163, 2023.

[38] Simone Scardapane, Michele Scarpiniti, Enzo Baccarelli, and Aurelio Uncini. Why should we add early exits to neural networks? *Cognitive Computation*, 12:954 – 966, 2020.

[39] Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Tran, Yi Tay, and Donald Metzler. Confident adaptive language modeling. In *NeurIPS*, 2022.

[40] Tal Schuster, Adam Fisch, Tommi S. Jaakkola, and Regina Barzilay. Consistent accelerated inference via confident adaptive transformers. In *EMNLP*, pages 4962–4979, 2021.

[41] Roy Schwartz, Gabriel Stanovsky, Swabha Swayamdipta, Jesse Dodge, and Noah A. Smith. The right tool for the job: Matching model and instance complexities. In *ACL*, pages 6640–6651, 2020.

[42] Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with pointer-generator networks. In *ACL*, pages 1073–1083, 2017.

[43] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *ArXiv*, abs/1909.08053, 2019.

[44] Qingyi Si, Tong Wang, Zheng Lin, Xu Zhang, Yanan Cao, and Weiping Wang. An empirical study of instruction-tuning large language models in chinese. *ArXiv*, abs/2310.07328, 2023.

[45] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Anand Korthikanti, Elton Zhang, Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi, Yuxiong He, Michael Houston, Saurabh Tiwary, and Bryan Catanzaro. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *ArXiv*, abs/2201.11990, 2022.

[46] Shengkun Tang, Yaqing Wang, Caiwen Ding, Yi Liang, Y. Li, and Dongkuan Xu. Deediff: Dynamic uncertainty-aware early exiting for accelerating diffusion model generation. *ArXiv*, abs/2309.17074, 2023.

[47] Surat Teerapittayanon, Bradley McDanel, and H. T. Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *ICPR*, pages 2464–2469, 2016.

[48] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023.

[49] Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288, 2023.

[50] Neeraj Varshney, Agneet Chatterjee, Mihir Parmar, and Chitta Baral. Accelerating llama inference by enabling intermediate layer decoding via instruction tuning with lite. *ArXiv*, abs/2310.18581, 2023.

[51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017.

[52] Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. Deebert: Dynamic early exiting for accelerating BERT inference. In *ACL*, pages 2246–2251, 2020.

[53] Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. Berxit: Early exiting for BERT with better fine-tuning and extension to regression. In *EACL*, pages 91–104, 2021.

[54] Canwen Xu and Julian McAuley. A survey on dynamic neural networks for natural language processing. In *EACL*, pages 2370–2381, 2023.

[55] Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *NeurIPS*, pages 12360–12371, 2019.

[56] Linfeng Zhang, Chenglong Bao, and Kaisheng Ma. Self-distillation: Towards efficient and compact neural networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(8):4388–4403, 2022.

[57] Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian J. McAuley, Ke Xu, and Furu Wei. BERT loses patience: Fast and robust inference with early exit. In *NeurIPS*, 2020.